

**NASA Contractor Report 182034**  
**ICASE Report No. 90-20**

# ICASE

## **THE COST OF CONSERVATIVE SYNCHRONIZATION IN PARALLEL DISCRETE EVENT SIMULATIONS**

**David M. Nicol**

Contract No. NAS1-18605  
May 1990

Institute for Computer Applications in Science and Engineering  
NASA Langley Research Center  
Hampton, Virginia 23665-5225

Operated by the Universities Space Research Association

(NASA-CR-182034) THE COST OF CONSERVATIVE  
SYNCHRONIZATION IN PARALLEL DISCRETE EVENT  
SIMULATIONS Final Report (ICASE) 32 p

CSCD 09B

N90-23912

Unclass

63/61 0280815



National Aeronautics and  
Space Administration

**Langley Research Center**  
Hampton, Virginia 23665-5225



# The Cost of Conservative Synchronization in Parallel Discrete Event Simulations

*David M. Nicol*  
*Department of Computer Science \**  
*College of William and Mary*

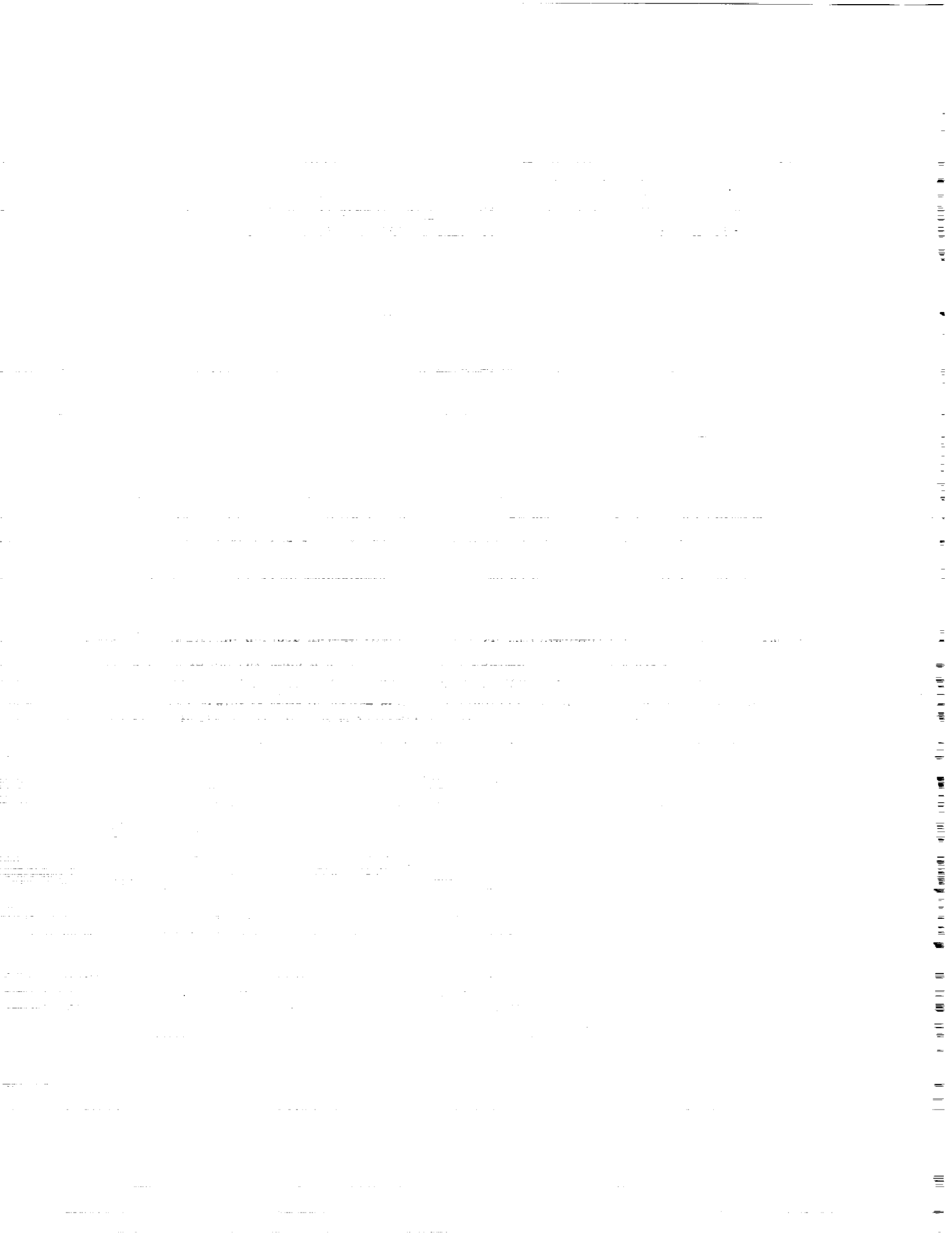
May 7, 1990

## Abstract

This paper analytically studies the performance of a synchronous conservative parallel discrete-event simulation protocol. The class of simulation models considered are oriented around a physical domain, and possess a limited ability to predict future behavior. Using a stochastic model we show that as the volume of simulation activity in the model increases relative to a fixed architecture, the complexity of the average per-event overhead due to synchronization, event list manipulation, lookahead calculations, and processor idle time approaches the complexity of the average per-event overhead of a serial simulation. The method is therefore within a constant factor of optimal. Our analysis demonstrates that on large problems—those for which parallel processing is ideally suited—there is often enough parallel workload so that processors are not usually idle. We also demonstrate the viability of the method empirically, showing how good performance is achieved on large problems using a thirty-two node Intel iPSC/2 distributed memory multiprocessor.

---

\*Supported in part by the Virginia Center for Innovative Technology, by NASA grants NAG-1-060 and NAS1-18605, and NSF grant ASC 8819393



# 1 Introduction

The problem of parallelizing discrete-event simulations has received a great deal of attention in the last several years. Simulations pose unique synchronization constraints due to their underlying sense of time. When the simulation model can be simultaneously changed by different processors, actions by one processor can affect actions by another. One must not simulate any element of the model too far ahead of any other in simulation time, to avoid the risk of having its logical past affected. Alternately, one must be prepared to fix the logical past of any element determined to have been simulated too far.

Two schools of thought have emerged concerning synchronization. The *conservative* school [5], [13], [23], [24] employs methods which prevent any processor from simulating beyond a point at which another processor might affect it. These synchronization points need to be re-established periodically to allow the simulation to progress. Early efforts focussed on finding protocols which were either free from deadlock, or which detected and corrected deadlock [17]. The *optimistic* school [7] allows a processor to simulate as far forward in time as it wants, without regard for the risk of having its simulation past affected. If its past is changed (due to interaction with a processor farther behind in simulation time) it must then be able to “rollback” in time at least that far, and must cancel any erroneous actions it has taken in its false future.

Conservative protocols are sometimes faulted for leaving processors idle, due to overly pessimistic synchronization assumptions. It is almost always true that individual *model elements* are blocked because of pessimistic synchronization; the conclusion that *processors* tend to be blocked requires the assumption that all model elements assigned to a processor tend to be blocked simultaneously, or that each processor has only one model element. The latter assumption pervades many performance studies, and is unrealistic for fine-grained simulation models executed on coarser grained multiprocessors. Intuition suggests that if there are many model elements assigned to each processor, then it is unlikely that *all* model elements on a processor will be blocked. Given sufficient workload, a properly designed conservative method should not leave processors idle, because there is so much work to do. While some model elements are blocked due to synchronization concerns, other elements, with high probability, are not.

It is natural to ask how much performance degradation due to blocking a conservative method suffers. We answer that question, by analyzing a simple conservative synchronization method. The method assumes the ability to pre-sample activity duration times[20], and assumes that any queueing discipline used is non-preemptive. The protocol itself is quite simple. As applied to a queueing network it works as follows. First, whenever a job enters service, the queue to which the job will be routed is immediately notified of that arrival (sometime in the future), and the receiving queue computes a service time for the new arrival. These two actions constitute *lookahead*, a concept which is key to the protocol’s success. Now imagine that all events with time-stamps less than  $t$  have already been processed and that the processors are globally synchronized. For each queue we determine the time-stamp of the

next job it would route (excluding one in service) if no further arrivals occur at that queue. The processors cooperatively compute the minimum such time, say  $\delta(t)$ . We will show that all further messages to be sent in the simulation have time-stamps at least as large as  $\delta(t)$ . Consequently a processor may evaluate, in parallel with all other processors, all of its events with time-stamps less than  $\delta(t)$ . Having done so, the processors synchronize globally, and repeat the process. The interval  $[t, \delta(t))$  is called a *window*, and  $\delta(t) - t$  is its *width*.

We analyze the performance of the protocol by first deriving an approximated lower bound on the equilibrium mean window width. We then multiply this width by the equilibrium rate at which the simulation generates events. The resulting product is an approximated lower bound on the the average number of events that are processed within a window. We then identify conditions under which the average number of events processed in a window increases without bound as the system simulation event generation rate increases. Next we analyze the synchronization, idle time, lookahead calculation, and event-list overheads of the protocol as a function of  $T$ , events in the system at a time. The average overhead per processed event is shown to be  $O(f(T))$ , where  $f(T)$  is the complexity of the average per-event overhead in a optimized serial simulation. Therefore the protocol's asymptotic performance (as  $T \rightarrow \infty$ ) is within a constant factor of optimal. Finally, we demonstrate the viability of the protocol empirically. A parallel simulation system based on the protocol has been implemented on a thirty-two node Intel iPSC/2 distributed memory multiprocessor[2]. Processor efficiencies in the range of 60% - 90% are reported for several different large simulation models.

It is important to remember that our analysis concerns average case performance based on a general stochastic model. Specific problem examples can be constructed to ensure that the protocol essentially executes serially, while another can execute many things in parallel. We believe that such examples are somewhat artificial and do not shed a great deal of light on how performance will behave over a wide range of problems. Our intention is to study the average case performance on a model of typical simulation problems.

This paper makes two basic contributions. One is to develop a new approach for the analysis of parallel discrete-event simulations. The second is a demonstration that many large simulation models having much concurrent activity can be effectively simulated in parallel using a simple conservative protocol.

This paper is organized as follows. §2 gives some background for this work. §3 describes the model of discrete-event simulations we use in our protocol description and analysis, and then introduces the protocol. §4 derives an approximated lower bound on the average number of events processed in a window. §5 determines the complexity of the average total overhead per event suffered by using the protocol. §6 reports on the performance of the protocol on several different simulation models. §7 gives our conclusions.

## 2 Background

Our protocol is similar to others recently proposed [1], [4], [13], [14], [28]. Unlike earlier asynchronous protocols, these synchronously move a window across simulation time, roughly as follows. Let *floor* be the lower edge of the window. This means that all events with time-stamps less than *floor* have already been processed. The processors then cooperatively determine the upper window edge, *ceiling*. This value is chosen in such a way that all events within  $[floor, ceiling)$  can safely be processed in parallel. *ceiling* becomes *floor* for the next window, and so on.

The major question with synchronous conservative protocols is whether windows small enough to prevent dependencies between window events admit enough such events to keep all the processors busy. Lubachevsky was the first to answer this question [14], by deriving a lower bound on the number of events processed within a window defined by his method. Using this bound and some assumptions concerning event density (in simulation time), he shows that the performance of his method scales up as the problem size and number of processors are simultaneously increased. However, his results are not quantitative, although they might have been so developed. Our analysis is different, in that we define a model from which event densities follow naturally, and we quantify the average number of events processed in a window. Ours is an average case analysis, while Lubachevsky's is a worst case analysis. Also, Lubachevsky's analysis hinges on the assumption of a non-zero minimal propagation delay, while ours does not. We *do* show that minimum service times can dramatically improve the average number of events processed each window.

The protocol we study is an application of the one described by Chandy and Sherman [4] to a more restricted problem domain. Like Lubachevsky's method, they require periodic global synchronization among processors. Each window their protocol computes the minimum time-stamp among all "conditional" events, and then processes all "unconditional" events with smaller time-stamps. In addition, their technique incorporates the conversion of "conditional" events into "unconditional" events, as a function of messages exchanged in the simulation. Such conversion is highly application dependent. The most important difference between our protocol and the general conditional-event approach lies in the specificity of our conversion of conditional events into unconditional events, in a way that requires little model-specific information. Furthermore, our protocol is stated within the context of a model closer to those used by simulation practitioners than is the model used to describe the conditional-event approach.

Our analysis of lookahead is related to that developed by Lin and Lazowska in [10], and by Wagner and Lazowska in [30]. Their work analyzes the ability of different queue types to predict future behavior, and focuses on lookahead at a single queue. Our analysis is of a much simpler lookahead scheme, but analyzed over the entire simulation. The protocol we describe can be easily adapted to accommodate these more complex techniques for computing lookahead. We have also analyzed a different class of simulations than the one studied here, on massively parallel architectures[19]. The sensitivity of performance to lookahead is

quantified, upper bounds on optimal and optimistic performance are derived, as is a lower bound on the performance of the same protocol we study in this paper.

Some analysis exists of the optimistic "Time Warp" method of synchronization. The earliest analyses concerned detailed stochastic models of two processor systems [9, 18]. These models include overhead costs and permit heterogeneous processors. Most other studies of Time Warp tend to assume negligible state-saving and rollback costs. For example, Lin and Lazowska have shown that if Time Warp has no state-saving or rollback costs, and if "correct" computations are never rolled back, then Time Warp achieves optimality [11]. This is intuitive, because Time Warp aggressively searches for the simulation's critical path—if it is able to do so without cost, its performance must be optimal. Other analyses highlight the fact that Time Warp can "guess right" while conservative methods must block. Lipton and Mizell have shown that there is a certain asymmetry between optimistic and conservative methods: while it is possible for an optimistic method to arbitrarily outperform a conservative method, the converse is not true [12]. Their analysis explicitly includes overhead costs. Madisetti, Walrand, and Messerschmitt [16] have developed a performance model which estimates the rate at which simulation time advances under an optimistic strategy such as Time Warp. They model the behavior of the system as a Markov chain, and include the cost of communication and of synchronization. Their analysis is exact for two processors, and approximate for a general number of processors. Lubachevsky, Schwartz, and Weiss use a sophisticated stochastic model to show how it is possible for Time Warp simulations to thrash in periods of "cascading rollbacks" [15].

### 3 Model and Protocol

We now describe our model of discrete-event simulations more formally, and define the synchronization protocol.

#### 3.1 Model Assumptions

Consider a domain containing  $S$  sites, where *activities* occur. An activity (e.g., service given to a job at a queue) begins, ends, and upon its completion enables (i.e. causes) other activities. These causations are reported to the appropriate sites by way of *completion* messages. Consequently, three distinct events are associated with each activity: **enable**, **begin**, and **complete**. The **enable** event for a given activity can be different from the **begin** event if the site imposes queueing. We permit a completion to cause more than one activity in order to include simulation problems such as Petri-nets, where a single transition firing may cause token arrivals at multiple Petri-net places. Thus, we assume that a **complete** event at site  $S_i$  causes an activity at each member of a random subset of other sites. All **enable** events caused by a completion have the same time-stamp as the completion. An activity is said to be *occurring* at time  $t$  if its associated **begin** event has a time-stamp no greater than  $t$ , and its **complete** event has time-stamp no less than  $t$ . Each site maintains



its own priority queue of events associated with activities enqueued or occurring at the site. Each site also maintains its own simulation *clock*, which records the time-stamp on the last event processed.

Under the assumptions of our model the **enable** and **complete** events are *unconditional*—once placed on the event list, no further activity in the simulation will change them. **begin** events may be conditional. For example, a **begin** event at time  $t$  might describe the future placement of a particular job into service at a queue at time  $t$ . If before  $t$  another job with higher priority arrives, that **begin** event may be removed from the event list.

Depending on the ability of the site, activities may occur there one at a time, or concurrently. We assume that either an unbounded number of activities may simultaneously occur at a site, or that only one activity may occur at a time. In the former case, we say the site has *infinite servers*. In the latter case, enabled activities may be enqueued before occurring. The delay in simulation time between when an activity begins and ends is called its *duration*. We assume that a duration is strictly positive, but do not assume a minimal duration. For the purposes of analysis we assume that the simulation model is ergodic, and that each duration time comes from a distribution composed by adding a nonnegative constant to an exponentially distributed random variable. Each site may have a unique distribution.

Our performance analysis rests on a number of assumptions about the simulation model which are exploited by the protocol.

1. We assume that once an activity begins, the causation of further activities cannot affect its completion time.
2. We assume that the simulation state change due to an activity completion is very local—the state change is implied by knowledge of which activity completed, which activities are subsequently caused, and the time of the completion.
3. We assume that the activities caused by the completion of activity  $A_j$  can be reported to their respective sites *at the time that  $A_j$  begins*.
4. We assume that a lower bound on the duration of an activity can be determined at the time of the receipt of the **completion** message which causes the activity.

To illustrate these assumptions, consider a job  $J$  which at time  $s$  begins service at a non-preemptive queue  $Q_1$ , completes at time  $\hat{s}$ , and is routed to  $Q_2$ . Assumption 1 is satisfied by the nature of  $Q_1$ 's queueing discipline. Assumption 2 is satisfied because the change in model state due to this departure is completely characterized by knowledge of  $Q_1$ ,  $Q_2$ , and  $\hat{s}$ . Assumption 3 is satisfied if the service discipline at  $Q_1$  is non-preemptive and the routing is independent of the jobs enqueued at time  $\hat{s}$ : in the simulation we can report the arrival of  $J$  at time  $\hat{s}$  to  $Q_2$  concurrently with the entering of  $J$  into service at  $Q_1$ , at time  $s$ . By doing so, the processing required of  $J$ 's **completion** event at  $\hat{s}$  does not include reporting  $J$ 's departure, but may include the recording of statistics which depend on all simulation activity at  $Q_1$  (including arrivals) up to time  $\hat{s}$ . Assumption 4 is satisfied if  $J$ 's service time

at  $Q_2$  can be computed at the time that  $Q_1$  reports the arrival of  $J$  to  $Q_2$ . This is possible if the service time of every job at  $Q_2$  is drawn independently from the same probability distribution.

This model describes a large number of common simulation models, and is related to *event graphs* described in [26] and [27]. Many queueing networks are obviously captured. Logic networks are described, with activities corresponding to logical module evaluations. Here new activities are caused when a module output changes state. The movement and interaction of objects in a domain can also be captured. One assumes no queueing at sites, and models the passage of an object across some discrete region of the domain as an activity. *Lookahead* plays a major role in our synchronization method and its analysis. Lookahead exists and is exploited by assumptions 3 and 4 above.

Simulation workload is the event processing. This includes changing anticipated event times as a result of newly caused activities, in changing simulation state variables, and in gathering/recording statistics. We view event list management costs as inescapable overheads associated with the processing of events.

Our protocol does not require a minimal duration time for its correctness. However, performance is substantially enhanced if every duration time is bounded from below by  $D_{\min} > 0$ . Equivalently, we can introduce a minimal time  $D_{\min}$  delay between when an activity completes, and when activities it causes are enabled. We will use  $D_{\min}$  throughout our analysis, but may take it to be zero.

### 3.2 Protocol Definition

Next we define the synchronization protocol in terms of the model given in §3.1. Our only architectural assumptions are that the simulation model is executed on a multiprocessor having  $P$  processors; any processor can send a message (indirectly, if needed) to any other processor, and the processors can synchronize globally.

One important aspect of our protocol is the “pre-sending” of completion messages. Let  $A_j$  be some activity whose **begin** event has time-stamp  $s$ . Let  $\hat{s}$  be  $A_j$ ’s completion time. Under our protocol  $A_j$ ’s site must send completion messages to all sites where activities caused by  $A_j$ ’s completion will occur, *at the time  $A_j$  begins*. Observe that even though the simulation time at  $A_j$ ’s site is  $s$ , these completion messages are time-stamped with time  $\hat{s} > s$ . A site which receives such a notification inserts an **enable** event with time-stamp  $\hat{s}$  into its event list (a non-queueing site may directly insert a **begin** event with time  $\hat{s}$ ); it also selects a duration time (or a lower bound on it) for the newly caused activity.

Suppose the processors have globally synchronized, and let  $t$  be the minimum time-stamp among events at all sites. Each site  $S_i$  can determine a lower bound  $\delta_i(t)$  on the earliest completion time of any of its pending (i.e., as yet not begun) activities, assuming no further **enable** events are received. We call this the site’s *lookahead bound*. For example, consider a site  $S_i$  with queueing. There are three cases to consider.

**Case 1:**  $S_i$ ’s event list is void of **enable** events. In this case we define  $\delta_i(t) = \infty$ .

**Case 2:** *No activity is occurring at  $t$ , and  $S_i$ 's event list contains enable events.* Let  $u$  be the earliest enable time among these, and define  $\delta_i(t)$  to be the completion time of the activity enabled at  $u$ .

**Case 3:** *Some activity is occurring at  $t$ , and  $S_i$ 's event list contains enable events.* Define  $\delta_i(t)$  to be the completion time of the next enabled activity to receive service, assuming that no further enable events will be inserted into the event list.

If  $S_i$  has infinite servers, only two cases arise. If there are no **begin** events in  $S_i$ 's event list, then define  $\delta_i(t) = \infty$ . If there are **begin** events in  $S_i$ 's event list, define  $\delta_i(t)$  to be the minimum completion time among these.

Finally, define

$$\delta(t) = \min_{\text{all sites } S_i} \{\delta_i(t)\}.$$

The protocol is very simple. Define  $w_1 = 0$ , and proceed as follows.

1. Given  $w_n$ , the processors cooperatively determine  $\delta(w_n)$ .
2. Each site may be simulated in parallel with all others until the time of the event with least time-stamp at that site is as large as  $\delta(w_n)$ . The processing of any **begin** event in this interval must include pre-sending the associated completion messages.
3. Sites receive the messages sent during the processing of  $[w_n, \delta(w_n))$ , select duration times for the associated caused activities, and insert events into their event lists.
4.  $n = n + 1$ . Goto step 1.

The obvious question to ask of this protocol is whether the sites can safely process all events within a window. The protocol is safe if, once the window is established, no further messages with time-stamps less than the upper edge of the window will ever be sent. The following theorem establishes this fact.

**Theorem 3.1** *Let  $[w_n, \delta(w_n))$  be a window established by the protocol. Then every completion message sent during the processing of  $[w_n, \delta(w_n))$  has a time-stamp at least as large as  $\delta(w_n)$ .*

**Proof:** Completion messages are pre-sent by the processing of **begin** events. Let  $b_0, \dots, b_k$  be the times of all **begin** events in  $[w_n, \delta(w_n))$ , in increasing order. We use induction to show that for  $i = 0, \dots, k$ , the completion messages associated with the **begin** event at time  $b_i$  have time-stamps at least as large as  $\delta(w_n)$ . For the base case consider  $b_0$ , and let  $S_i$  be the associated site.  $S_i$  computes  $\delta_i(w_n)$  to be the minimum time-stamp on the next message it sends, *provided no further messages are received at  $S_i$* . By construction  $S_i$  will not receive any further messages with time-stamps less than  $b_0$ , therefore the decision to

begin the activity at  $b_0$  was correctly fore-seen during the computation of  $\delta_i(w_n)$ , implying that the completion time of the activity beginning at  $b_0$  is no smaller than  $\delta_i(w_n)$ , and hence is no smaller than  $\delta(w_n)$ . This establishes the base of the induction. For the induction step suppose that the completion times of the activities begun at times  $b_0, \dots, b_{j-1}$  are all no smaller than  $\delta(w_n)$ . Consider the activity begun at time  $b_j$ , and let  $S_i$  be its site. As a consequence of the induction hypothesis, during the processing of  $[w_n, \delta(w_n))$   $S_i$  cannot receive any messages with time-stamps less than  $b_j$ . Consequently, the decision to begin an activity at time  $b_j$  was correctly fore-seen during the computation of  $\delta_i(w_n)$ . The completion time of the activity beginning at  $b_j$  is thus no smaller than  $\delta_j(w_n)$ , and so is no smaller than  $\delta(w_n)$ . This completes the induction.

□

Under the assumption of non-zero duration times, it will always be true that  $w_n < \delta(w_n)$ . Consequently, simulation time advances each window (even if no events occur in the window), and deadlock never occurs.

### 3.3 Example

An example helps to illustrate the protocol's mechanism. Consider a system with sites  $S_1$  and  $S_2$ . Site  $S_1$  permits an unbounded number of activities to occur simultaneously, while site  $S_2$  imposes queueing. The system moves *objects* between sites. Duration times are random. When an object completes its duration it either disappears, moves to another (possibly the same) site, or splits into a number of objects that move.  $S_2$  uses Last-Come-First-Serve queueing.

Let  $w_n = 100$ , and imagine that objects  $O_1$  and  $O_2$  are present at  $S_1$ , with scheduled completion times of 100 and 103. Object  $O_3$  is in service at  $S_2$ , and will complete at time 101. Object  $O_4$  is enqueued at  $S_2$ , and will eventually receive 4 units of service.

The completion of  $O_1$  at time 100 sends  $O_1$  back to  $S_1$ , where it will receive another 8 units of service; the completion of  $O_2$  at time 103 sends  $O_2$  to  $S_2$  where it will eventually receive 6 units of service;  $O_2$ 's completion at time 103 also creates a new object  $O_5$  which is sent to  $S_1$ , where it receives 4 units of service. At site  $S_2$ ,  $O_3$  completes at time 101, and then remains at  $S_2$ , where it will receive another 5 units of service. Observe that the messages reporting the completions of  $O_1$ ,  $O_2$ , and  $O_3$  have already been sent, and the "next" durations of those objects have already been chosen.

This scenario is summarized in figure 1, along with the contents of  $S_1$  and  $S_2$ 's event lists as observed at time 100. The event lists reflect the practice of pre-sending object arrival notices.  $S_1$  determines its lookahead bound  $\delta_1(100)$  by finding the minimum completion time among all objects it knows will arrive at or after time 100.  $O_2$  arrives (again) at time 100 and completes at 108.  $O_5$  arrives at 103 and completes at 107, making  $\delta_1(100) = 107$ .  $S_2$  determines  $\delta_2(100)$  by identifying the next object to complete service that isn't already in service. Because  $S_2$  is LCFS, the arrival of  $O_3$  at time 101 causes  $O_3$  to receive service before  $O_4$ .  $\delta_2(100)$  is 106, so that  $\delta(100) = 106$ .  $S_1$  and  $S_2$  are thus free to simulate all events with

Objects at time 100						
<i>Object</i>	<i>Site</i>	<i>Arrives</i>	<i>Duration</i>	<i>Completes</i>	<i>Routed</i>	<i>Comments</i>
$O_1$	$S_1$	?	?	100	$S_1$	Occurring at time 100
$O_1$	$S_1$	100	8	108	?	Caused by completion of self
$O_2$	$S_1$	?	?	103	$S_2$	Occurring at time 100
$O_2$	$S_2$	103	6	?	?	Caused by completion of self
$O_5$	$S_1$	103	4	107	?	Caused by $O_2$ at $S_1$
$O_3$	$S_2$	?	?	101	$S_2$	Occurring at time 100
$O_3$	$S_2$	101	5	106	?	Highest priority activity at 101
$O_4$	$S_2$	?	4	?	?	Lower priority at 101

$\delta(100) = 106$

Event Lists

$S_1$ Event List at time 100		$S_2$ Event List at time 100	
<i>Event</i>	<i>Time</i>	<i>Event</i>	<i>Time</i>
$O_1$ completes	100	$O_3$ completes	101
$O_1$ arrives	100	$O_3$ arrives	101
$O_2$ completes	103	$O_2$ arrives	103
$O_5$ arrives	103		
4 events processed in [100, 106)		3 events processed in [100, 106)	

Figure 1: Example of Synchronous Protocol Operation

times no greater than 106, in parallel.  $S_1$  has four such events,  $S_2$  has three (or four, if the processing of the  $O_3$  arrival event at 101 creates a **begin** event at 101).

Arrival events (**enable** events) at  $S_1$  may also serve as **begin** events since no queueing is imposed. Each site's processing of arrival events includes the decision of where to route the object upon completion, and the generation of completion messages with the appropriate time-stamp.

## 4 Analysis of Protocol

Our performance analysis derives an approximated lower bound on the mean window width, then multiplies by the equilibrium event creation rate in order to bound the average number of events created per window. By flow balance this bounds the average number of events processed per window. We then consider the behavior of this average as a function of

simulation activity rate, and minimum duration time.

The analysis to follow uses results from the theory of stochastic order relations, and manipulates hazard rate functions. Readers unfamiliar with these tools should consult Ross [25]; the appendix quickly sketches the main ideas and results we use.

We are interested in the limiting value of the expected window width  $E[\delta(w_n) - w_n]$  as  $n \rightarrow \infty$ , supposing that the limit exists. As we will see, a window's width is comprised of the minimum of a number of complicated random variables. Complications arise both due to randomness in the model (e.g., random selection of sites where activities are caused following a completion), and due to dependence of the random variables' distributions on the past activity in the simulation. Our approach is to bound the mean window width from below with the mean minimum of much simpler, and stochastically smaller, random variables. The stochastically smaller variables are constructed by considering hazard rate functions. This is a useful analytic trick which exploits the fact that the hazard rate function for the minimum of a group of independent random variables is just the sum of their individual hazard rate functions.

One step in the bounding argument is intuitive, but not rigorously justified. Therefore one can only rigorously call our results approximate.

The analysis uses a slightly more formal model than we have yet described. The duration time distribution for site  $S_i$  is taken to be  $D_i + \exp\{\mu_i\}$ , where  $D_i \geq 0$  is constant and  $\exp\{\mu_i\}$  is exponential with mean  $\mu_i = 1/\lambda_i$ . We let  $D_{\min}$  be the minimum  $D_i$  value among all sites. The discussion of random variables, means, and hazard rates all concern the stochastic portion of the duration times.

Our bounds depend on the manner in which a completing activity causes activities elsewhere. To more precisely describe these effects, for every site  $S_i$  let  $Reach(S_i)$  be the set of all sites where activities caused by a completion at  $S_i$  can occur. For convenience we assume that the activities caused by a single completion are all at different sites. Activity  $A_j$  completing at  $S_i$  randomly chooses a subset  $B_j \subseteq Reach(S_i)$ , and causes one activity at each site in  $B_j$ . We assume  $B_j$  is chosen independently of the duration values of the caused activities. The distribution governing this choice is particular to  $S_i$ ;  $p(B, i)$  denotes the probability that  $B \subseteq Reach(S_i)$  is the selected set.

Let  $A_j$  be an activity occurring at site  $S_i$ , and let  $B$  be the set of sites with activities caused by  $A_j$ . We will be interested in the rate at which the first activity completes, among all those caused by  $A_j$ . Towards this end, we focus on the stochastic portion of these activity durations. The "rate" of the minimum stochastic portion is just  $\lambda_B = \sum_{S_j \in B} \lambda_j$  (see §A.2). The expected rate (with respect to the distribution of  $B$ ) is defined by

$$\begin{aligned} \psi_i &= \sum_{B \subseteq Reach(S_i)} p(B, i) \left( \sum_{S_j \in B} \lambda_j \right) \\ &= \sum_{S_j \in Reach(S_i)} \Pr\{\text{completion at } S_i \text{ causes an activity at } S_j\} \lambda_j. \end{aligned} \quad (1)$$

Pathological analytic difficulties are avoided by assuming that the simulation model al-

ways has at least one activity occurring that causes other activities. This can be ensured, for example, by adding a “clock” site that does nothing but process a single, periodically self-causing activity.

Let  $\mathcal{A}(w_n)$  be the random set of activities occurring at time  $w_n$ . During most of the analysis to follow we will condition on knowing that  $\mathcal{A}(w_n)$  is some fixed subset  $V$ . The conditioning is undone later when we take an expectation with respect to the distribution of  $\mathcal{A}(w_n)$ .

The discussion to follow focuses on activities. To facilitate precise reference to the site where a given activity occurs, we often use the notation  $S_{s(j)}$  to describe the site at which activity  $A_j$  occurs.

Consider the construction of the  $n$ th window. To compute  $\delta(w_n)$  we examine each site to determine the time of the next message it will send, in the absence of receiving any further messages. For a non-queueing site  $S_i$  this is the minimum completion time among all activities with **begin** events in  $S_i$ ’s event list. For each such activity there is another which caused it, and which is occurring at time  $w_n$ . If  $S_i$  is a queueing site, the activity  $A_j$  whose completion defines  $S_i$ ’s lookahead bound is either enqueued waiting for the completion of an occurring activity at  $S_i$ , or has its **enable** event sometime in the future. In the latter case we know there must be another site with an activity which is occurring at time  $w_n$ , and which causes  $A_j$ . Therefore, every activity whose completion defines some lookahead bound can be associated with an activity occurring at  $w_n$ . Conversely, for every  $A_j \in V$  we can associate a set of sites  $C_j$  with activities caused by  $A_j$ , such that the completion time of each activity  $A_k \in C_j$  equals  $\delta_{s(k)}(w_n)$ .  $C_j$  is obviously a subset of  $B_j$ , the set of all sites with activities caused by  $A_j$ , so that the minimum completion time among all activities caused by  $A_j$  at sites in  $B_j$  is no larger than the minimum taken over  $C_j$ .

We will want to distinguish queueing sites from non-queueing sites. We therefore define the indicator coefficient  $\gamma_i$  to have value 1 if site  $S_i$  is a queueing site, and to have value 0 if not.

For every  $A_j \in V$  let  $R_j(w_n)$  denote the *residual* duration time of  $A_j$ —the difference between  $A_j$ ’s completion time and  $w_n$ . For every  $A_j \in V$  at a queueing site  $S_{s(j)}$  define  $N_j(w_n)$  to be  $\infty$  if  $\delta_{s(j)} = \infty$ , otherwise it is the duration of the enqueued activity whose completion time is  $\delta_{s(j)}(w_n)$ . Let  $E_j(w_n)$  be the enabling time of the activity defining  $N_j(w_n)$ . Observe that this activity is sensitive to  $w_n$ : if  $A_j$  was occurring at time  $w_{n-1}$  it is possible for a higher priority activity to be enabled between times  $w_{n-1}$  and  $w_n$ , so that the activities defining  $N_j(w_{n-1})$  and  $N_j(w_n)$  may be different.

We define  $N_j(w_n) = \infty$  if  $S_{s(j)}$  is a non-queueing site. Regardless of whether  $S_{s(j)}$  is a queueing or non-queueing site we may say that the completion rate of  $N_j(w_n)$ ’s stochastic portion is  $\gamma_{s(j)}\lambda_{s(j)}$ .

Again, let  $B_j$  be the set of sites with activities caused by  $A_j$ ; for each  $S_k \in B_j$  let  $D_k + Y_{j,k}$  be the duration of the activity at  $S_k$  caused by  $A_j$ . We define  $A_j$ ’s lookahead bound to be the minimum completion time among (i) the activities caused by  $A_j$ , (ii) the next activity to complete at  $S_{s(j)}$  if  $S_{s(j)}$  is non-queueing and receives no further **enable** events.  $A_j$ ’s

lookahead bound as measured at time  $w_n$  may be written as

$$K_j(w_n) = w_n + R_j(w_n) + \min\{\max\{0, E_j(w_n) - R_j(w_n)\} + N_j(w_n), \min_{S_k \in C_j} \{D_k + Y_{j,k}\}\}.$$

$\delta(w_n)$  is the minimum lookahead bound among all activities  $A_j \in V$ . We may therefore write

$$\begin{aligned} E[\delta(w_n) - w_n \mid \mathcal{A}(w_n) = V] &= E[\min_{A_j \in V} \{K_j(w_n)\}] \\ &\geq E[\min_{A_j \in V} \{R_j(w_n) + \min\{N_j(w_n), \min_{S_k \in B_j} \{D_k + Y_{j,k}\}\}\}] \end{aligned} \quad (2)$$

The expectation above is complicated by its dependence on the history of the synchronization behavior up to time  $w_n$ . For example, suppose that activity  $A_j$  began in the  $(n - b_j)$ th window, for some  $b_j > 0$ . The distribution of  $K_j(w_n)$  must be conditioned on the event  $\mathcal{G}_j(w_n)$  that  $K_j(w_{n-c}) \geq w_{n-c+1}$  for all  $1 \leq c < b_j$ . Since  $K_j(w_n)$  is largely comprised of random variables that also comprise  $K_j(w_{n-c})$  for each  $c$ , conditioning on  $\mathcal{G}_j(w_n)$  makes each  $K_j(w_n)$  probabilistically larger than it would be if each component random variable had its original, unconditional distribution. The starting point for our bound is to build a stochastically smaller replacement for each  $K_j(w_n)$  by replacing each of  $K_j(w_n)$ 's components with a pristine unconditional random variable with the appropriate distribution.

We construct an “unconditioned” lookahead variable for each  $A_j$  as follows. Randomly choose a subset  $\mathcal{U}_{s(j)} \subseteq \text{Reach}(S_{s(j)})$  in accordance with the probability distribution  $\{p(B, s(j))\}$ , and independently choose a duration time  $D_k + X_{j,k}$  for each  $S_k \in \mathcal{U}_{s(j)}$ .  $D_k + X_{j,k}$  will replace the actual corresponding duration time  $D_k + Y_{j,k}$ . Randomly and independently choose some value  $D_{s(j)} + W_{j,s(j)}$  from  $S_{s(j)}$ 's duration time distribution. If  $S_{s(j)}$  is a non-queueing site we take  $W_{j,s(j)} = \infty$ .  $D_{s(j)} + W_{j,s(j)}$  will replace the actual  $N_j(w_n)$ . Let  $Z_{j,s(j)}$  be an independent exponential having the distribution of the stochastic portion of  $S_{s(j)}$ 's duration time.  $Z_{j,s(j)}$  will replace  $R_j(w_n)$ ; note that the residual of an  $S_{s(j)}$  duration time is always as large as the residual of the duration time's stochastic portion.

The event  $\mathcal{G}_j(w_n)$  gives us information that  $K_j(w_n)$  is probabilistically larger than it would be if its components had their original distributions. Therefore, intuition suggests that the following inequality is true

$$\lim_{n \rightarrow \infty} E[\delta(w_n) - w_n] \geq \lim_{n \rightarrow \infty} E[\min_{A_j \in \mathcal{A}(w_n)} \{Z_{j,s(j)} + \min\{D_{s(j)} + W_{j,s(j)}, \min_{S_k \in \mathcal{U}_{s(j)}} \{D_k + X_{j,k}\}\}\}]. \quad (3)$$

Note that the expectations involved in this assumption are not conditioned on  $\mathcal{A}(w_n) = V$ , and that we only require the inequality to hold in the limit of  $n \rightarrow \infty$ . It seems exceedingly difficult to formally establish this bound. Our analysis therefore proceeds by assuming its validity.

**Assumption 4.1** *Inequality (3) is true.*



We continue the analysis by placing stochastic lower bounds on variables comprising the conditional (on  $\mathcal{A}(w_n) = V$ ) expectation

$$E\left[\min_{A_j \in V} \{Z_{j,s(j)} + \min\{D_{s(j)} + W_{j,s(j)}, \min_{S_k \in \mathcal{U}_{s(j)}} \{D_k + X_{j,k}\}\}\}\right]. \quad (4)$$

As a first step we note that

$$\min_{S_k \in \mathcal{U}_{s(j)}} \{D_k + X_{j,k}\} \geq \min_{S_k \in \mathcal{U}_{s(j)}} \{X_{j,k}\} + D_{\min}. \quad (5)$$

Next we put a stochastic lower bound on  $\min\{X_{j,k} | S_k \in \mathcal{U}_{s(j)}\}$ . This random variable is complicated by the fact that  $\mathcal{U}_{s(j)}$  is a random set. For any *given* set  $\mathcal{U}_{s(j)} = B_i$ , the minimum of  $|B_i|$  exponentials is itself an exponential, with rate  $\lambda_{B_i} = \sum_{S_k \in B_i} \lambda_k$ . Consequently  $\min\{X_{j,k} | S_k \in \mathcal{U}_{s(j)}\}$  is a probabilistic mixture of exponentials—with probability  $p(B_i, s(j))$  it is an exponential with rate  $\lambda_{B_i}$ . Without loss of generality we may enumerate all subsets  $B_i \subseteq \text{Reach}(S_{s(j)})$  in such a way that  $\lambda_{B_i} \leq \lambda_{B_j}$  whenever  $i < j$ . Given this ordering, Lemma A.1 establishes that an exponential  $T_{j,s(j)}$  whose rate is the “expected” rate  $\psi_{s(j)} = \sum_{B_i} p(B_i, s(j)) \lambda_{B_i}$  (see expression (1)) is stochastically smaller than the minimum:

$$\min_{S_k \in \mathcal{U}_{s(j)}} \{X_{j,k}\} \geq_{st} T_{j,s(j)}.$$

Applying inequalities (5) and (11) we determine that

$$\begin{aligned} & \min\{D_{s(j)} + W_{j,s(j)}, \min_{S_k \in \mathcal{U}_{s(j)}} \{D_k + X_{j,k}\}\} \\ & \geq \min\{D_{s(j)} + W_{j,s(j)}, \min_{S_k \in \mathcal{U}_{s(j)}} \{X_{j,k}\} + D_{\min}\} \\ & \geq_{st} \min\{D_{\min} + W_{j,s(j)}, T_{j,s(j)} + D_{\min}\} \\ & = \min\{W_{j,s(j)}, T_{j,s(j)}\} + D_{\min}. \end{aligned} \quad (6)$$

Since  $W_{j,s(j)}$  and  $T_{j,s(j)}$  are both exponential, their minimum is also exponential and has rate  $\gamma_{s(j)} \lambda_{s(j)} + \psi_{s(j)}$  (recall that  $\gamma_{s(j)} = 0$  and  $W_{j,s(j)} = \infty$  if  $S_{s(j)}$  is a non-queueing site). Let  $U_{j,s(j)}$  be an exponential with rate  $\gamma_{s(j)} \lambda_{s(j)} + \psi_{s(j)}$ . Inequality (6) holds for every  $A_j \in V$ ; furthermore, the lookahead random variable constructed for each  $A_j$  is independent of all others. Since the addition and min operators are increasing it follows from (11) that the expectation in (4) is bounded from below:

$$\begin{aligned} & E\left[\min_{A_j \in V} \{Z_{j,s(j)} + \min\{D_{s(j)} + W_{j,s(j)}, \min_{S_k \in \mathcal{U}_{s(j)}} \{D_k + X_{j,k}\}\}\}\right] \\ & \geq E\left[\min_{A_j \in V} \{Z_{j,s(j)} + U_{j,s(j)}\}\right] + D_{\min}. \end{aligned} \quad (7)$$

We remove the conditioning on  $V$  by taking the expectation with respect to  $\mathcal{A}(w_n)$ ,

$$\begin{aligned} & E\left[\min_{A_j \in \mathcal{A}(w_n)} \{Z_{j,s(j)} + \min\{D_{s(j)} + W_{j,s(j)}, \min_{S_k \in \mathcal{U}_{s(j)}} \{D_k + X_{j,k}\}\}\}\right] \\ & \geq E\left[\min_{A_j \in \mathcal{A}(w_n)} \{Z_{j,s(j)} + U_{j,s(j)}\}\right] + D_{\min}. \end{aligned}$$

If  $\mathcal{A}$  has the limiting distribution of  $\mathcal{A}(w_n)$  as  $n \rightarrow \infty$  (supposing it exists), then

$$\begin{aligned} \lim_{n \rightarrow \infty} E[ \min_{A_j \in \mathcal{A}(w_n)} \{ Z_{j,s(j)} + \min\{ D_{s(j)} + W_{j,s(j)}, \min_{S_k \in \mathcal{U}_{s(j)}} \{ D_k + X_{j,k} \} \} \} ] \\ \geq E[ \min_{A_j \in \mathcal{A}} \{ Z_{j,s(j)} + U_{j,s(j)} \} ] + D_{\min}. \end{aligned} \quad (8)$$

Our next task is to deal with the randomness of the set  $\mathcal{A}$ .

Let the collection of site servers in the domain be enumerated as  $V_1, V_2, \dots$ , and define  $v(j)$  to be the index of  $V_j$ 's site. For each  $i = 1, 2, \dots, S$  and  $j = 1, 2, \dots$  let

$$\omega_i = \lim_{n \rightarrow \infty} E[\text{number of site } S_i \text{ activities occurring at time } w_n],$$

$$\rho_j = \lim_{n \rightarrow \infty} \Pr\{\text{at time } w_n \text{ an activity is occurring at } V_j\},$$

and observe that

$$\omega_i = \sum_{V_j, v(j)=i} \rho_j,$$

assuming that the expectations and limits exist. It is not obvious that  $\omega_i$  should be identical to the equilibrium expected number of activities occurring at  $S_i$ ; intuitively one expects it to be close, because the number of windows in which a given activity is found occurring is roughly proportional to the duration of the activity.

The expectation on the right-hand-side of (8) is taken with respect to a distribution of random sets of activities found occurring at a window edge. One can equivalently view it as an expectation taken with respect to a random set of servers found busy at a window edge. Inequality (7) suggests we associate two exponentials with each server  $V_j$ :  $Z_j$  and  $U_j$  (here binding  $j$  to the server rather than to the activity). There is a one-to-one correspondence between a random subset of servers, and a random subset  $H \subseteq \{(Z_1, U_1), (Z_2, U_2), \dots\}$ .

Lemma A.2 was developed to deal with the situation at hand. Following its statement we define

$$\begin{aligned} \Lambda &= \sum_{j=1}^{\infty} \Pr\{(Z_j, U_j) \in H\} \lambda_{v(j)} (\gamma_{v(j)} \lambda_{v(j)} + \psi_{v(j)}) \\ &= \sum_{j=1}^{\infty} \rho_j \lambda_{v(j)} (\gamma_{v(j)} \lambda_{v(j)} + \psi_{v(j)}) \\ &= \sum_{i=1}^S \omega_i \lambda_i (\gamma_i \lambda_i + \psi_i). \end{aligned} \quad (9)$$

The lemma's conclusion is that

$$E[ \min_{(Z_j, U_j) \in H} \{ Z_j + U_j \} ] \geq \sqrt{\frac{\pi}{2\Lambda}}.$$

The left-hand-side of this inequality is identical to the right-hand-side of (8), except for the inclusion of  $D_{\min}$ . Assuming the validity of assumption 4.1 we may conclude that

$$\lim_{n \rightarrow \infty} E[\delta(w_n) - w_n] \geq D_{\min} + \sqrt{\frac{\pi}{2\Lambda}}. \quad (10)$$

In order to determine the average number of events processed per window we need to consider the rate at which events are generated by the simulation. Let  $\tilde{\omega}_i$  be the equilibrium mean number of activities occurring at  $S_i$ . There are two events associated with each activity at a non-queueing site, **begin** and **complete**. Adding **enable**, there are three events associated with an activity at a queueing site. We therefore define the variable  $\epsilon_i$  to be 2 or 3 depending on whether  $S_i$  is non-queueing or queueing, respectively.

An activity's duration at  $S_i$  has mean  $\lambda_i^{(d)} = 1/(D_i + \mu_i)$ , so that the equilibrium event creation rate is  $\Lambda_{Sys} = \sum_{i=1}^S \epsilon_i \tilde{\omega}_i \lambda_i^{(d)}$ . By flow balance this is also the equilibrium event completion rate. We can therefore multiply this rate times the lower bound on the mean window width to bound the mean number of events processed in a window.

**Theorem 4.2** *Let*

$$\Lambda_{Sys} = \sum_{i=1}^S \epsilon_i \tilde{\omega}_i \lambda_i^{(d)}$$

*be the system event creation rate, and let*

$$\Lambda = \sum_{i=1}^S \omega_i \lambda_i (\gamma_i \lambda_i + \psi_i).$$

*Then if assumption 4.1 is valid, the average number of events processed per window is at least*

$$\Lambda_{Sys} \left( D_{\min} + \sqrt{\frac{\pi}{2\Lambda}} \right).$$

□

This theorem demonstrates how an existing minimal service time accelerates performance. Given constant  $D_{\min} > 0$ , the bound increases at least linearly as the total simulation event rate increases. However, good performance is also possible when  $D_{\min} = 0$ , as we will see.

The value of  $\Lambda$  is defined in terms of  $\omega_i$ . We have no immediate cause for believing that  $\omega_i = \tilde{\omega}_i$ ; nor is it clear that the two quantities should be widely different. It seems reasonable then to take  $\omega_i \approx \tilde{\omega}_i$  as a first approximation. Doing so permits us to analytically estimate  $\Lambda$  in some simple cases, and quantify the bound given by Theorem 4.2.

As pointed out by Wagner and Lazowska [30], interconnection topology plays an important role in determining the performance one achieves with a queueing system. Network bottlenecks limit the volume of simulation activity. This is reflected in Theorem 4.2. For example, in a network where each site has one server,  $\omega_i$  is approximately the server utilization. A bottleneck site will have a very high utilization while those at other sites are comparatively low. After a point, adding jobs to the network does not appreciably increase the sum of server utilizations, hence the overall event rate does not appreciably increase. For the same reason simulated queueing systems are constrained even if the throughput at each site is equal. The overall system event rate is maximized when all site utilizations are one. After a point, to increase simulation activity one needs to increase the size of the network.

We can approximate the bound in Theorem 4.2 in some simple cases. Consider a model where objects move throughout the domain. An object resides at a site for a fixed time  $D_{\min}$  plus an exponential time with mean  $1/\lambda$ , and then moves to any other site, chosen uniformly at random. Equilibrium flow balance equations are easily solved in this situation. Working through the details with  $K$  objects and  $D_{\min} = 0$ , one discovers that at least  $\sqrt{\pi K/2}$  events are processed per window, on average. A relevant point is that the inter-site communication topology is that of a fully connected graph. Such topologies are generally taken to be extremely taxing on conservative synchronization methods, because the “next” event at a site can come from anywhere. Nevertheless, a significant amount of work is performed each window, at least when  $K$  is large. Figure 2 plots the analytically bounded and empirically measured average number of events processed per window, as a function of  $\log_2 K$ . The empirical measurements represent the sample mean of ten long simulation runs. There was very little variance between these runs. Figure 2 shows that if thousands of objects are in the model, hundreds of events are processed each window. Since parallel processing techniques are used primarily when serial processing times are too slow (or memories are too small), we see that this result applies directly to situations of practical interest—large simulation models on medium scale parallel architectures.

Performance is greatly enhanced when  $D_{\min} > 0$ . Figure 3 plots measurements of the number of events per window for small models, having only 256 and 1024 objects. The same measurement methodology as was described for Figure 2 is used here. The analytic bound is not displayed, being indistinguishable from the measured performance when plotted on the graph.  $D_{\min}$  is varied between 0 and  $\mu = 1/\lambda$ , so that  $D_{\min}/\mu$  varies between 0 and 1. We see that if a model has minimal duration times we can expect many more events per window than if not. Note that the protocol does not need to know  $D_{\min}$ , as it is already part of the pre-sampled duration times. Dramatic performance improvement as one’s ability to “look ahead” increases has also been observed by Fujimoto [6].

Our confidence in the conclusions of Theorem 4.2 is increased by the fact that the approximated lower bound did uniformly fall below measured performance. Similar results have been observed when comparing the measured and bounded performance on less homogeneous simulation models.

## 5 The Cost of Conservative Synchronization

Next we consider the overheads involved in implementing this conservative protocol. First we identify conditions under which the average number of events processed per window will grow without bound as the system event creation rate grows without bound. Then we show that as the number of events processed per window grows, our method’s per-event overhead due to synchronization, processor idle time, lookahead calculation, and event list manipulation becomes within a constant factor of average the per-event overhead of performing the simulation serially.

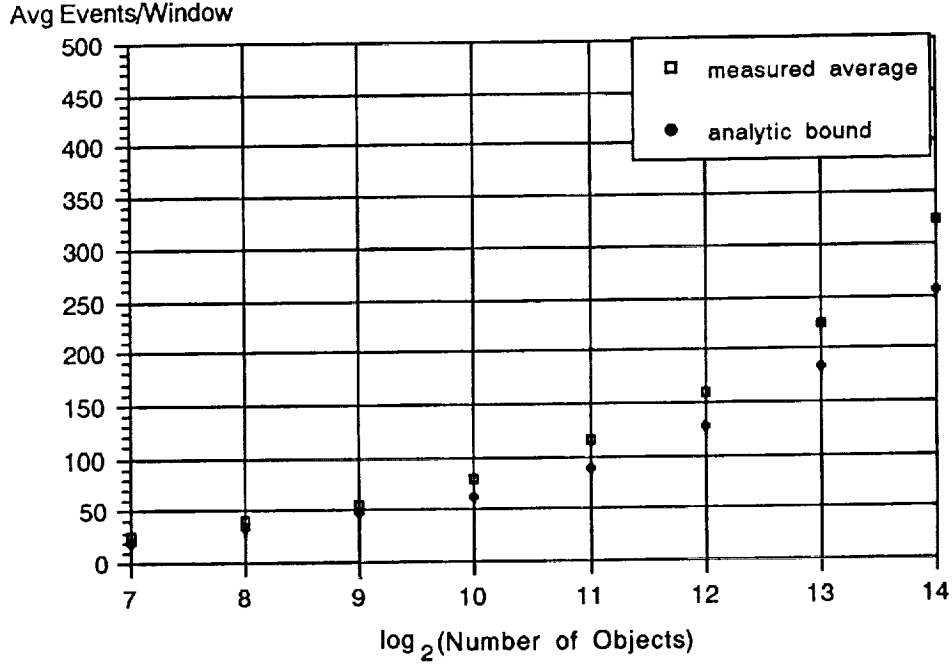


Figure 2: Average events per window, as function of number of objects.  $D_{\min} = 0$ ; durations are homogeneous exponentials; no queueing; routing is uniformly random.

One way of increasing the system event creation rate  $\Lambda_{sys}$  is to increase the “size” of the model. For example, we increase the size of the moving objects simulation described earlier by increasing the number of objects in the domain. We may also increase the number of sites, although in this case it is not necessary. Theorem 4.2 shows how the average number of events processed each window may increase as  $\Lambda_{sys}$  increases. Clearly, if  $D_{\min} > 0$  then at least  $\Lambda_{sys} D_{\min}$  events are processed each window on average. It is also possible for the average number of events to increase without bound as  $\Lambda_{sys}$  increases even when  $D_{\min} = 0$ . For example, suppose there is a value  $\alpha$  such that as the size of the simulation model is increased the following bound is true for all sites  $S_i$ :

$$\frac{\omega_i \lambda_i (\lambda_i + \psi_i)}{\tilde{\omega}_i \lambda_i^{(d)}} \leq \alpha.$$

This condition is a formal statement that as the model size grows  $\psi_i$  can’t get too large relative to  $\lambda_i$ , and that any difference between  $\omega_i$  and  $\tilde{\omega}_i$  doesn’t get too large. The first condition will be satisfied if there exists  $\lambda_{\max}$  and  $R_{\max}$  such that as the model size grows,  $\lambda_i \leq \lambda_{\max}$  and  $|Reach(S_i)| \leq R_{\max}$ , for all  $i$ . The second condition ought to be satisfied

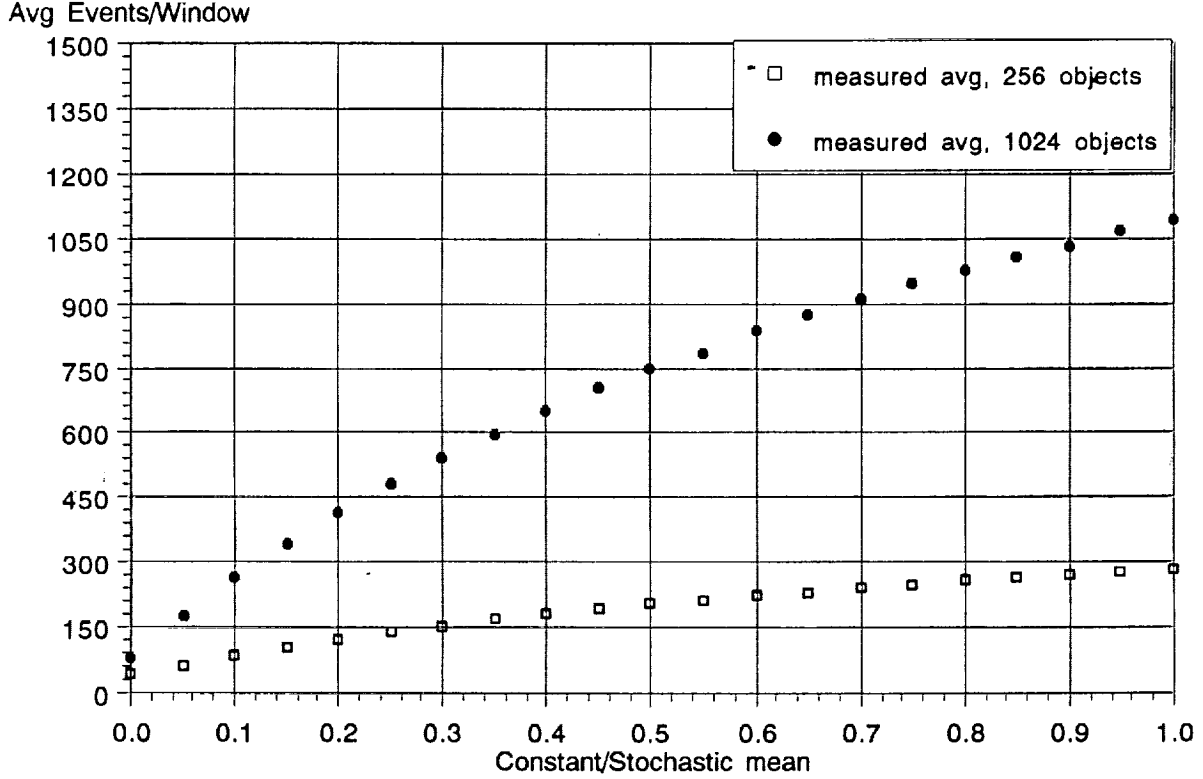


Figure 3: Average events per window when  $D_{\min} > 0$ . Performance plotted as function of  $D_{\min}/\mu$ , for 256 and 1024 objects.

if our intuition that  $\omega_i \approx \tilde{\omega}_i$  is correct. If the bound above holds, then as the model size grows the inequality  $\Lambda \leq \alpha \Lambda_{Sys}$  will always hold. It follows that at least  $\sqrt{\pi \Lambda_{Sys}/(2\alpha)}$  events are processed each window on average, a number that grows without bound as  $\Lambda_{Sys}$  grows without bound.

As a point of comparison, we assume that a serial implementation uses the best known event list management algorithm. If there are  $T$  total events in the system on average, we let  $f(T)$  be the average complexity of an optimized serial event list algorithm. For example, there is some evidence that a “calendar-queue” implementation has an average  $O(1)$  complexity (i.e.,  $f(T) = 1$ ) on the hold model [3]. A number of other event list algorithms exhibiting  $\Omega(\log T)$  average complexity are also commonly used [8]. We assume that the serial event list algorithm permits the deletion of a non-minimal element without affecting the overall average complexity. This assumption is satisfied by the calendar queue implementation.

We make the reasonable assumption that as the simulation model size is increased,  $T$

grows at least as rapidly as  $S$ , i.e.,  $T = \Omega(S)$ .

Now consider a parallel simulation that uses our protocol. The requirement that completion notices be pre-sent may increase the average total number of events in event lists at a time. This represents a factor of two increase, at most. In the complexity analysis to follow we need not explicitly concern ourselves with this constant factor increase.

One overhead suffered in the parallel simulation is the cost of computing  $\delta_i(w_n)$  for each site  $S_i$ . This value changes only when an event is inserted or deleted at  $S_i$ . A queueing site can recompute this value with  $O(1)$  cost whenever its event list changes. A non-queueing site can organize the completion times of its pending activities in a priority queue using whatever event list algorithm is employed by the optimized serial implementation. The minimum value in the priority queue defines  $\delta_i$ . The priority queue is modified only when the site's event list is modified, at cost  $O(f(T))$ . A processor can organize the  $\delta_i$  values from each of its sites into priority queue, enabling it to determine the minimum on-processor  $\delta_i$  value at least as quickly as the optimized serial implementation finds its minimal element. Maintenance of this priority queue costs  $O(f(S))$  on average for each processed event. Once each processor has determined its locally minimum  $\delta_i$  value, all processors may cooperatively compute the global minimum in  $C_{Synch}$  time. Note that our assumption that  $P$  is fixed permits us to ascribe a worst-case constant cost to this operation.

Another overhead is processor idle time. The protocol is punctuated with global synchronizations, between which the processors execute in parallel. A processor with little workload will spend a long time waiting for more heavily loaded processors to reach the synchronization barrier. Suppose there are  $W$  events to process in a window. For the purposes of analysis, assume that each event may be mapped to any processor, with equal probability<sup>1</sup>. Then the number of events assigned to a processor is a binomial  $B(W, 1/P)$  random variable. The collection of workload random variables are not independent however, as we know they must sum to  $W$ . However, it isn't difficult to construct a coupling[25] argument to show that the expected maximum workload of this system must be smaller than the expected maximum workload of a system where each processor has an independent  $B(W, 1/P)$  workload. The binomial distribution has an increasing hazard rate function [25](p.280); it is therefore stochastically less variable than an exponential with the same mean [25](p.273), and hence the expected maximum of  $P$  independent exponential random variables with mean  $W/P$  is at least as large as the expected maximum of  $P$  independent  $B(W, 1/P)$  random variables. The expected maximum of the exponentials is approximately  $(W/P) \ln(W/P)$ . Assuming each event takes the same amount of time to process, the average fraction of time a processor is left idle is no greater than

$$1 - \frac{(W/P)}{(W/P) \ln P} = 1 - \frac{1}{\ln P}.$$

This implies that the average overhead cost per event due to processor idleness is  $O(1)$ . This

---

<sup>1</sup>This can't rigorously be true, since events at the same site are evaluated on the same processor. It is a reasonable approximation when  $W$  is large compared with  $P$ .

analysis is actually quite pessimistic. Much better load balance can be achieved through use of mapping techniques such as scatter decomposition[22]. Also, the bound above is insensitive to increasing volume of workload, whereas in practice the proportion of idle time tends to decrease as the volume of workload increases.

The complexity of the average per-event overhead due to event list manipulation, lookahead calculations, processor idle time, and global synchronization is

$$\frac{W(O(f(T)) + O(f(S)) + O(1)) + C_{Synch}}{W} = O(f(T)).$$

Relative to the serial simulation, performance must then be within a constant factor of optimal, at least if inter-*LP* communication costs are ignored (we have already accounted for the communication costs that are specific to the protocol). Inter-*LP* communication costs are dependent on the simulation model and its mapping, and are dependent on the architecture. It is possible for communication costs to overwhelm performance, even if our protocol finds a great deal of parallel workload. However these costs are inherent to the model, and would be suffered under any synchronization protocol.

## 6 Empirical Results

We used the protocol analyzed in this paper in a parallel discrete-event simulation testbed implemented on the Intel iPSC/2 distributed memory multiprocessor[2]. The testbed, YAWNS (Yet Another Windowing Network Simulator)[21], is designed to permit rapid development of simulation models, by providing a framework within which all synchronization and inter-processor communication activity is automated, and hidden from the user. YAWNS uses a computational paradigm where the simulation model is decomposed into communicating *Logical Processes (LPs)*. *LP*'s interact by passing messages. A site in our analytic model plays the role of an *LP*.

The simulation modeler must provide the testbed with three routines for each *LP* (the *LP*'s may share these routines). One routine processes messages, typically inserting an event into the *LP*'s event list as a result. This routine is responsible for choosing a duration time for the enabled activity. Another routine processes events. Messages to other *LP*'s may be generated as a result of calling this routine; these messages correspond to the completion messages described in the analytic model. The third routine is called to obtain the lookahead value required of an *LP*. YAWNS demands that the simulation modeler know about the protocol only to the extent that inter-*LP* messages are pre-sent, and an *LP* must be able to determine a lower bound on the time of the next message it sends.

It is always important to use the best possible event list algorithm for an *LP*. YAWNS provides a linearly-linked list algorithm for use when the number of events in an *LP*'s list is small, and a splay-tree algorithm for large lists.

We report on the performance achieved by four diverse applications: the moving objects simulation described earlier, a logic network, the game of Life, and a timed Petri net model.



All measurements reported are taken from a thirty-two processor machine. Each simulation model was run long enough to generate several millions of events. The execution time was typically a minute or two, once the problem was loaded and running. Much longer runs were also performed, but no appreciable difference in performance statistics were observed.

The measured performance supports our analysis, and actually becomes quite good on large problems. The metric we use to gauge performance is average processor utilization, measured as the fraction of time a processor spends doing work that would be performed in a serial implementation of the simulation, using the same paradigm. Time spent in computing lookahead, synchronization, interprocessor communication, and idle time are explicitly counted as overhead, and do not appear in the utilization figure. One can translate such efficiencies into “speedup” figures by multiplying by the number of processors used, provided the resulting numbers are properly interpreted. The speedup so computed is relative to a serial version that uses the same paradigm (and code) of communicating *LPs* as is used in the parallel version. This is not an unreasonable paradigm for a general purpose serial simulation system, but is not likely to be the paradigm of choice for a serial version that is highly optimized for the given application. In our experience (and depending on the application), the communicating *LP* paradigm is a factor of 1.5 to 2 times slower than an optimized serial version. The usual comparison of serial running time to parallel running time is impossible to directly obtain, as the largest models we simulate are too large for a single processor’s memory. We will see that on the largest problems the average processor utilization ranges from 60% – 90%.

## 6.1 Moving Objects

The sites are connected in a hypercube topology. In each model there are exactly as many objects as there are sites. Each object resides at a site for a time constructed by adding 0.25 to an exponential with mean 1. We increase the size of the problem by simultaneously increasing the number of objects and the number of sites. We may therefore describe the size of the system by the dimension of the underlying hypercube. Pre-sent completion times and lookahead values are computed exactly as described for non-queueing sites in this paper. Average processor utilization  $\rho$  as a function of hypercube dimension is given below. Many simulation runs were performed, the variance in the timing numbers is quite small.

Dim	8	9	10	11	12	13	14
$\rho$	21%	28%	34%	46%	54%	60%	62%

## 6.2 Logic Network

To ensure that we simulated networks with high concurrency we constructed “random” logic networks having the topology of a butterfly interconnection network. The last stage wraps around to feed the first. Each gate was randomly assigned to be an AND, OR, or XOR

function and was given a randomly chosen gate delay time of 1,2, or 3 time units. Each gate was modeled as an *LP*. The eventual output of a gate whose inputs have changed can be computed at the time the inputs change, hence gate state changes can be pre-sent. A gate is like a non-queueing site; its lookahead is computed to be the gate delay plus the minimum time of the next input change. The size of network can be described by the dimension of a column of gates. For example, a network of dimension 6 has 6 columns, each composed of  $2^6$  gates. Observed performance is given below.

Dim	5	6	7	8	9	10	11
$\rho$	24%	32%	43%	52%	59%	66%	70%

### 6.3 Conway's Game of Life

Initial random configurations were chosen so that the probability of a cell being alive is at step 0 is 0.2. Each cell is modeled as an *LP*. A cell is evaluated at step  $n$  only if one of its neighbors (or itself) is alive at step  $n - 1$ . It is straightforward to pre-send "new state" messages; lookahead consists of one step time. The problem size is increased by increasing the size of the board. Again, we can easily describe problem size in terms of dimension. A  $2^j \times 2^j$  board will be said to have dimension  $j$ .

Dim	3	4	5	6	7	8
$\rho$	12%	16%	35%	54%	69%	77%

Larger problems than a  $256 \times 256$  board will often exhaust the available dynamic memory in some processor, after some period of execution. This points out one of consequences of internally buffering all messages until the window's workload is completed.

### 6.4 Timed Petri Nets

Consider a timed Petri net model of a multiprocessor system organized with a mesh communication topology. The net models a system where a processor iteratively receives a message from each of its NEWS (North, East, West, South) neighbors, performs a computation, and sends a result to each NEWS neighbor. The net models a flow control policy that prevents a processor from sending a message to a neighbor until the last message it sent to that neighbor is consumed. An *LP* consists of the network for one processor, a network containing approximately thirty places and ten transitions. Nearly all transitions have a unit time delay associated with them. Transitions modeling the processor execution time have 200 units of delay.

This Petri net model does not satisfy exactly the assumptions we've made concerning simulation model behavior. The main difference is that a token arriving to an *LP* does not trigger a single *LP* activity with a single duration time. The response of the *LP* is liable

to be much more complex. Nevertheless, the basic synchronization protocol works. Tokens from an enabled transition are always pre-sent (regardless of whether they are sent to places within the *LP*); to compute lookahead, an *LP* adds the minimum delay among all transitions that send tokens to other *LPs* to the least-time token arrival event in the *LP*'s event list.

The grid size for the simulated system can be described in terms of dimension in the same way as was the Game of Life.

Dim	3	4	5	6
$\rho$	35%	62%	84%	94%

The comparatively better performance of this problem can be attributed to its better ratio of computation costs to *LP*-message costs.

## 7 Conclusions

We have analyzed a simple conservative synchronization protocol for parallel discrete-event simulation. The protocol presumes that one can pre-sample activity duration times (or bound those times from below), that the immediate effects of simulation model state changes are very local, and that all queueing disciplines are non-preemptive. The protocol essentially slides a window across simulation time; the window is defined so that processors can evaluate all their window events in parallel. We construct an approximated lower bound on the average number of events processed per window. The bound depends on the topology and activity rates of the heterogeneous simulation domain. The performance analysis shows that a great deal of workload can be performed in parallel, if there is a great deal of concurrent activity in the simulation model. Non-zero minimal activity durations are shown to greatly improve performance. We show that the asymptotic time complexity of the average total overhead (synchronization, lookahead calculations, processor idle time, event list manipulation) per event is that of an optimized serial simulation. Assuming that the complexity of the communication cost per event is no greater than the overhead of an event in a serial implementation, the protocol's performance is within a constant factor of optimal. The region of problems where the method does well is precisely the region where parallel processing is most effectively applied—problems too large to run serially. The method is verified by implementation on a distributed memory multiprocessor. Good performance is observed on a variety of problems.

## A Appendix

In this appendix we describe the tools used in our analysis, and develop some key results.

## A.1 Stochastic Dominance

Our analysis relies on the theory of stochastic dominance. The definitions and results we cite are taken from Ross [25], chapter 8.

Random variable  $X$  is said to be *stochastically larger* than random variable  $Y$  if for all  $t$

$$\Pr\{X > t\} \geq \Pr\{Y > t\}.$$

We then write  $X \geq_{st} Y$ , or  $Y \leq_{st} X$ . An equivalent definition is that

$$E[g(X)] \geq E[g(Y)] \quad \text{for all increasing functions } g.$$

In particular,  $E[X] \geq E[Y]$ . If  $X_1, \dots, X_n$  are independent random variables and  $Y_1, \dots, Y_n$  are independent random variables such that  $X_i \geq_{st} Y_i$  for all  $i$ , then for all increasing functions  $g$ ,

$$g(X_1, \dots, X_n) \geq_{st} g(Y_1, \dots, Y_n). \quad (11)$$

## A.2 Hazard Rate Functions

If  $X$  is a nonnegative continuous random variable, it has a *hazard rate function*, also known as a *failure rate function*. Let  $f(t)$  be  $X$ 's density function, and let  $\bar{F}(t) = \Pr\{X > t\}$ . Then  $X$ 's hazard rate function is defined to be

$$\lambda(t) = \frac{f(t)}{\bar{F}(t)}.$$

If  $X$  is exponential, then  $\lambda(t)$  is identically the exponential's rate parameter.

We rely on the following results concerning hazard rate functions.

- If  $\lambda_X(t)$  and  $\lambda_Y(t)$  are hazard rate functions for  $X$  and  $Y$ , and  $\lambda_X(t) \leq \lambda_Y(t)$  for all  $t$ , then  $X \geq_{st} Y$ .
- If  $X_1, \dots, X_n$  are independent random variables with hazard rate functions  $\lambda_1(t), \dots, \lambda_n(t)$ , then the hazard rate function for  $\min\{X_1, \dots, X_n\}$  is simply  $\sum_{i=1}^n \lambda_i(t)$ .
- If  $X$  has hazard rate function  $\lambda(t)$ , then for any  $t$  and  $s$ ,  $s \leq t$ ,

$$\Pr\{X > t | X > s\} = \exp\left\{-\int_s^t \lambda(u) du\right\}.$$

This also shows (taking  $s = 0$ ) that the hazard rate function uniquely defines a distribution.

### A.3 Important Bounds

We now establish some important bounds used in this paper.

Random variables constructed by randomly choosing one of a set of random variables are called *mixtures*. The following lemma bounds the hazard rate of a certain class of mixtures.

**Lemma A.1** *Let  $X_1, X_2, \dots$ , be independent random variables with hazard rate functions  $\lambda_1(t), \lambda_2(t), \dots$ , and suppose that these functions are ordered:*

$$\lambda_i(t) \leq \lambda_{i+1}(t) \quad \text{for all } i = 1, 2, \dots, \text{ and all } t \geq 0.$$

*Let  $p_1, p_2, \dots$  be probabilities such that  $\sum_{i=1}^{\infty} p_i = 1$ , and consider the random variable  $M$  constructed by randomly selecting some  $X_i$ , with probability  $p_i$ . Let  $\lambda_M(t)$  be  $M$ 's hazard rate function. Then for all  $t \geq 0$*

$$\lambda_M(t) \leq \sum_{i=1}^{\infty} p_i \lambda_i(t).$$

**Proof:** <sup>2</sup> Let  $f_i(t)$  and  $\bar{F}_i(t)$  be the density and cumulative distribution functions for  $X_i$ . Then  $\lambda_i(t) = f_i(t)/\bar{F}_i(t)$ , and

$$\lambda_M(t) = \frac{\sum_{i=1}^{\infty} p_i f_i(t)}{\sum_{i=1}^{\infty} p_i \bar{F}_i(t)}.$$

The desired conclusion will follow if

$$\sum_{i=1}^{\infty} p_i f_i(t) \leq \left( \sum_{i=1}^{\infty} p_i \bar{F}_i(t) \right) \left( \sum_{i=1}^{\infty} p_i \lambda_i(t) \right)$$

for all  $t$ . Let  $Y = i$  with probability  $p_i$  and let  $h(Y, t) = \lambda_Y(t)$  and  $g(Y, t) = -\bar{F}_Y(t)$ . Then for every fixed  $t$ ,  $h$  and  $g$  are increasing in  $Y$ . Application of Proposition 7.1.5 of [25](p. 227) yields

$$E[h(Y, t)]E[g(Y, t)] \leq E[h(Y, t)g(Y, t)],$$

or equivalently,

$$\sum_{i=1}^{\infty} p_i f_i(t) \leq \left( \sum_{i=1}^{\infty} p_i \bar{F}_i(t) \right) \left( \sum_{i=1}^{\infty} p_i \lambda_i(t) \right).$$

As this holds for every  $t \geq 0$ , the lemma's conclusion follows.

□

We now develop a lower bound on the expected minimum of a random number of variables, each variable being the sum of two exponentials.

---

<sup>2</sup>This elegant proof was suggested by an anonymous referee, replacing a far more complicated proof of our own.

**Lemma A.2** Let  $\mathcal{S} = \{(Z_1, U_1), (Z_2, U_2), \dots\}$  be a countable set where  $Z_i$  is exponential with rate  $\lambda_i$ , and  $U_i$  is exponential with rate  $\psi_i$ . Let all these random variables be independent. Let  $B_1, B_2, \dots$  be the set of finite subsets of  $\mathcal{S}$ . Let  $B$  be a random set constructed by choosing  $B_i$  with probability  $p_i$ . Let

$$\Lambda = \sum_{i=1}^{\infty} \Pr\{(Z_i, U_i) \in B\} \lambda_i \psi_i.$$

Then

$$E\left[\min_{(Z_i, U_i) \in B} \{Z_i + U_i\}\right] \geq \sqrt{\frac{\pi}{2\Lambda}}.$$

**Proof:** Consider the hazard rate function  $\gamma_i(t)$  for  $Z_i + U_i$ . This random variable is the lifetime of a serial two-stage system where the first stage lasts for time  $Z_i$ , and the second lasts for time  $U_i$ .  $\gamma_i(t)$  is the instantaneous probability density associated with the system dying at time  $t$ , given that it has survived up to time  $t$ . Now if  $Z_i > t$ , the system cannot fail at  $t$ , whence  $\gamma_i(t) = 0$ . If  $Z_i \leq t$ , then the hazard rate is simply that of  $U_i$ :  $\psi_i$ . Note that this observation relies on the memoryless property of the exponential. We may therefore write

$$\begin{aligned} \gamma_i(t) &= (1 - \Pr\{Z_i > t \mid Z_i + U_i > t\})\psi_i \\ &\leq (1 - \Pr\{Z_i > t\})\psi_i \\ &= (1 - \exp\{-t\lambda_i\})\psi_i. \end{aligned}$$

One can show that the left-hand-side of this inequality is equivalent to the more usual (and complicated) derivation of the hazard rate function for the sum of two exponentials [29](p. 126). The function on the right-hand-side is concave in  $t$ , and is hence dominated everywhere by the line tangent to it at  $t = 0$ :  $\tau_i(t) = t\lambda_i\psi_i$ . A random variable  $V_i$  with hazard rate function  $\tau_i(t)$  satisfies  $V_i \leq_{st} Z_i + U_i$ .

Let  $B_j$  be any finite subset of  $\mathcal{S}$ . By (11) and the observations above we may conclude that

$$E\left[\min_{(Z_i, U_i) \in B_j} \{Z_i + U_i\}\right] \geq E\left[\min_{(Z_i, U_i) \in B_j} \{V_i\}\right].$$

We now focus on the right-hand-side of this inequality. The hazard rate function for  $M_j = \min\{V_i \mid (Z_i, U_i) \in B_j\}$  is simply

$$\lambda_{B_j}(t) = t \cdot \left( \sum_{(Z_i, U_i) \in B_j} \lambda_i \psi_i \right).$$

Without loss of generality we may enumerate the finite subsets of  $\mathcal{S}$  in such a way that if  $i < j$ , then  $\lambda_{B_i}(t) \leq \lambda_{B_j}(t)$  for all  $t$ . Let  $M$  be a mixture of  $\{M_1, M_2, \dots\}$ , where  $M_j$  is chosen with probability  $p_j$ ; let  $\lambda_M(t)$  be  $M$ 's hazard rate function. By Lemma A.1 we can

## REFERENCES

bound  $\lambda_M(t)$  from above by  $\lambda_Y(t)$ , defined by

$$\begin{aligned}\lambda_M(t) &\leq \sum_{i=1}^{\infty} p_i \lambda_{B_i} \\ &= \sum_{i=1}^{\infty} \Pr\{(Z_i, U_i) \in B\} t \lambda_i \psi_i = \lambda_Y(t).\end{aligned}$$

Let  $Y$  be a random variable with hazard rate function  $\lambda_Y(t)$ . Using the correspondence between hazard rate functions and probability distributions (see §A.2), we have

$$\begin{aligned}\Pr\left\{\min_{(Z_i, U_i) \in B} \{V_i\} > t\right\} &\geq \Pr\{Y > t\} \\ &= \exp\left\{-\int_0^t \lambda_Y(s) ds\right\} \\ &= \exp\left\{-\sum_{i=1}^{\infty} \Pr\{(Z_i, U_i) \in B\} t^2 \lambda_i \psi_i / 2\right\} \\ &= \exp\{-\Lambda t^2 / 2\}.\end{aligned}$$

Now

$$\begin{aligned}E\left[\min_{(Z_i, U_i) \in B} \{Z_i + U_i\}\right] &= \int_0^{\infty} \Pr\{\min\{Z_i + U_i | (Z_i, U_i) \in B\} > t\} dt \\ &\geq \int_0^{\infty} \exp\{-\Lambda t^2 / 2\} dt \\ &= (1/\sqrt{\Lambda}) \int_0^{\infty} \exp\{-s^2 / 2\} ds \quad \text{by defining } s = t\sqrt{\Lambda} \\ &= \sqrt{\frac{\pi}{2\Lambda}}.\end{aligned}$$

□

## References

- [1] AYANI, R. A parallel simulation scheme based on distances between objects. In *Distributed Simulation 1989*, SCS Simulation Series, 1989, pp. 113–118.
- [2] BOMANS, L., AND ROOSE, D. Benchmarking the iPSC/2 hypercube multiprocessor. *Concurrency: Practice and Experience* 1, 1 (Sept. 1989), 3–18.
- [3] BROWN, R. Calendar queues: a fast  $O(1)$  priority queue implementation for the simulation event set problem. *Communications of the ACM* 31, 10 (Oct. 1988), 1220–1227.
- [4] CHANDY, K., AND SHERMAN, R. The conditional event approach to distributed simulation. In *Distributed Simulation 1989*, SCS Simulation Series, 1989, pp. 93–99.

## REFERENCES

- [5] CHANDY, K. M., AND MISRA, J. Distributed simulation: a case study in design and verification of distributed programs. *IEEE Trans. on Software Engineering* 5, 5 (September 1979), 440-452.
- [6] FUJIMOTO, R. M. Performance measurements of distributed simulation strategies. In *Distributed Simulation 1988*, SCS Simulation Series, 1988, pp. 14-20.
- [7] JEFFERSON, D. R. Virtual time. *ACM Trans. on Programming Languages and Systems* 7, 3 (1985), 404-425.
- [8] JONES, D. An empirical comparison of priority-queue and event-set implementations. *CACM* 29, 4 (April 1986), 300-311.
- [9] LAVENBERG, S., AND MUNTZ, R. Performance analysis of a rollback method for distributed simulation. In *Performance '83*, Elsevier Science Pub. B. V. (North Holland), 1983, pp. 117-132.
- [10] LIN, Y., AND LAZOWSKA, E. *Exploiting Lookahead in Parallel Simulation*. Tech. Rep., Dept. of Computer Science, Univ. of Washington, 1989.
- [11] LIN, Y., AND LAZOWSKA, E. Optimality considerations for "Time Warp" parallel simulation. In *Distributed Simulation 1990*, SCS Simulation Series, 1990, pp. 29-34.
- [12] LIPTON, R., AND MIZELL, D. Time Warp vs. Chandy-Misra: a worst-case comparison. In *Distributed Simulation 1990*, SCS Simulation Series, 1990, pp. 137-143.
- [13] LUBACHEVSKY, B. Efficient distributed event-driven simulations of multiple-loop networks. *Communications of the ACM* 32, 1 (1989), 111-123.
- [14] LUBACHEVSKY, B. Scalability of the bounded lag distributed discrete-event simulation. In *Distributed Simulation 1989*, SCS Simulation Series, 1989, pp. 100-107.
- [15] LUBACHEVSKY, B., SHWARTZ, A., AND WEISS, A. Rollback sometimes works...if filtered. In *Proceedings of the 1989 Winter Simulation Conference* (Washington, D.C., 1989), pp. 630-639.
- [16] MADASETTI, V., WALRAND, J., AND MESSERSCHMITT, D. Synchronization in message-passing computers: models, algorithms, and analysis. In *Distributed Simulation 1990*, SCS Simulation Series, 1990, pp. 35-48.
- [17] MISRA, J. Distributed discrete-event simulation. *Computing Surveys* 18, 1 (March 1986), 39-60.
- [18] MITRA, D., AND MITRANI, I. Analysis and optimum performance of two message-passing parallel processors synchronized by rollback. In *Performance '84*, Elsevier Science Pub. B. V. (North Holland), 1984, pp. 35-50.



## REFERENCES

- [19] NICOL, D. Analysis of synchronization in massively parallel discrete-event simulations. In *Proceedings of the 1990 ACM PPoPP Symposium* (Seattle, March 1990), pp. 89–98.
- [20] NICOL, D. Parallel discrete-event simulation of FCFS stochastic queueing networks. *SIGPLAN Notices* 23, 9 (September 1988), 124–137.
- [21] NICOL, D., MICHEAL, C., AND INOUE, P. Efficient aggregation of multiple LP's in distributed memory parallel simulations. In *Proceedings of the 1989 Winter Simulation Conference* (Washington, D.C., December 1989), pp. 680–685.
- [22] NICOL, D., AND SALTZ, J. An analysis of scatter decomposition. *IEEE Trans. on Computers* (1990). To appear.
- [23] PEACOCK, J. K., MANNING, E., AND WONG, J. W. Synchronization of distributed simulation using broadcast algorithms. *Computer Networks* 4 (1980), 3–10.
- [24] REYNOLDS, JR., P. A shared resource algorithm for distributed simulation. In *Proceedings of the Ninth Annual International Computer Architecture Conference* (Austin, Texas, April 1982), pp. 259–266.
- [25] ROSS, H. *Stochastic Processes*. Wiley, New York, 1983.
- [26] SARGENT, R. Event graph modelling for simulation with application to flexible manufacturing systems. *Management Science* 34, 10 (October 1988), 1231–1251.
- [27] SCHRUBEN, L. Simulation modeling with event-graphs. *Communications of the ACM* 26 (1983), 957–963.
- [28] SOKOL, L., BRISCOE, D., AND WIELAND, A. MTW: a strategy for scheduling discrete simulation events for concurrent execution. In *Distributed Simulation 1988*, SCS Simulation Series, 1988, pp. 34–42.
- [29] TRIVEDI, K. *Probability and Statistics, with Reliability, Queueing, and Computer Science Applications*. Prentice-Hall, Englewood Cliffs, NJ, 1982.
- [30] WAGNER, D., AND LAZOWSKA, E. Parallel simulation of queueing networks: limitations and potentials. In *Proceedings of the 1989 SIGMETRICS Conference* (Berkeley, CA, 1989), pp. 146–155.

# Report Documentation Page

1. Report No. NASA CR-182034 ICASE Report No. 90-20		2. Government Accession No.		3. Recipient's Catalog No.	
4. Title and Subtitle THE COST OF CONSERVATIVE SYNCHRONIZATION IN PARALLEL DISCRETE EVENT SIMULATIONS				5. Report Date May 1990	
				6. Performing Organization Code	
7. Author(s)  David M. Nicol				8. Performing Organization Report No. 90-20	
				10. Work Unit No. 505-90-21-01	
9. Performing Organization Name and Address Institute for Computer Applications in Science and Engineering Mail Stop 132C, NASA Langley Research Center Hampton, VA 23665-5225				11. Contract or Grant No. NAS1-18605	
				13. Type of Report and Period Covered Contractor Report	
12. Sponsoring Agency Name and Address National Aeronautics and Space Administration Langley Research Center Hampton, VA 23665-5225				14. Sponsoring Agency Code	
15. Supplementary Notes Langley Technical Monitor: Richard W. Barnwell  Submitted to the Journal of the ACM  <u>Final Report</u>					
16. Abstract  This paper analytically studies the performance of a synchronous conservative parallel discrete-event simulation protocol. The class of simulation models considered are oriented around a physical domain, and possess a limited ability to predict future behavior. Using a stochastic model we show that as the volume of simulation activity in the model increases relative to a fixed architecture, the complexity of the average per-event overhead due to synchronization, event list manipulation, lookahead calculations, and processor idle time approaches the complexity of the average per-event overhead of a serial simulation. The method is therefore within a constant factor of optimal. Our analysis demonstrates that on large problems---those for which parallel processing is ideally suited---there is often enough parallel workload so that processors are not usually idle. We also demonstrate the viability of the method empirically, showing how good performance is achieved on large problems using a thirty-two node Intel iPSC/2 distributed memory multiprocessor.					
17. Key Words (Suggested by Author(s))  discrete-event simulation, parallel processing, performance analysis			18. Distribution Statement  61 - Computer Programming and Software 66 - Systems Analysis  Unclassified - Unlimited		
19. Security Classif. (of this report) Unclassified		20. Security Classif. (of this page) Unclassified		21. No. of pages 31	
				22. Price A03	